



# Singular Zeros of Polynomial Systems

Angelos Mantzaflaris, Bernard Mourrain

## ► To cite this version:

Angelos Mantzaflaris, Bernard Mourrain. Singular Zeros of Polynomial Systems. SAGA – Advances in ShApes, Geometry, and Algebra. Geometry and Computing, vol 10., pp.77-103, 2014, 10.1007/978-3-319-08635-4\_5 . hal-02958889

**HAL Id: hal-02958889**

**<https://inria.hal.science/hal-02958889>**

Submitted on 6 Oct 2020

**HAL** is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

---

# Singular Zeros of Polynomial Systems

Angelos Mantzaflaris and Bernard Mourrain

GALAAD, INRIA Méditerranée  
BP 93, 06902 Sophia-Antipolis, France  
`[FirstName.LastName]@inria.fr`

**Summary.** Singular zeros of systems of polynomial equations constitute a bottleneck when it comes to computing, since several methods relying on the regularity of the *Jacobian matrix* of the system do not apply when the latter has a non-trivial kernel. Therefore they require special treatment. The algebraic information regarding an isolated singularity can be captured by a finite, local basis of differentials expressing the *multiplicity structure* of the point.

In the present article, we review some available algebraic techniques for extracting this information from a polynomial ideal. The algorithms for extracting the, so called, *dual basis* of the singularity are based on matrix-kernel computations, which can be carried out numerically, starting from an approximation of the zero in question.

The next step after obtaining the multiplicity structure is to *deflate* the root, that is, construct a new system in which the singularity is eliminated. Having a deflated system allows to refine the solution fast and to high accuracy, since the Jacobian matrix is regular and all the usual machinery, e.g. Newton's method or existence and unicity criteria may be applied. Standard verification methods, based e.g. on interval arithmetic and a fixed point theorem, can then be employed to certify that there exists a unique perturbed system with a singular root in the domain.

## 1 Introduction

A main challenge in algebraic and geometric computing is singular point identification and treatment. Such problems naturally occur when computing the topology of implicit curves or surfaces [1], the intersection of parametric surfaces in geometric modeling. When algebraic representations are used, this reduces to solving polynomial systems. Several approaches are available: algebraic techniques such as Gröbner bases or border bases, resultants, subdivision algorithms [16], [19], homotopies, and so on. At the end of the day, a numerical approximation or a box of isolation is usually computed to identify every real root of the polynomial system. But we often need to improve the numerical approximation of the roots. Numerical methods such as Newton's iteration

can be used to improve the quality of the approximation, provided that we have a simple root. In the presence of a multiple root, the difficulties are significantly increasing. The numerical approximation can be of very bad quality, and the methods used to compute this approximation are converging slowly (or not converging). The situation in practical problems, as encountered in CAGD for instance, is even worse, since the coefficients of the input equations are known with some incertitude. Computing multiple roots and root multiplicities of approximate polynomial systems is an ill-posed problem, since changing slightly the coefficients may transform a multiple root into a cluster of simple roots (or even make it disappear).

For instance Newton’s method converges only linearly to such a point, if it converges at all [6]. Also, certification tests for the existence of roots on a domain do not directly treat these cases. On the other hand, computing the local multiplicity structure around a singularity breaks down to stable linear algebra methods, which can be run approximately. One can use this local structure to deflate the root, and thus restore super-linear convergence of Newton iteration, or use standard verification techniques to certify a singular root of the original system. In case of inexact coefficients, known up to a certain tolerance, an exact singular root no longer exists. Nevertheless, a well chosen symbolic perturbation, combined with deflation, allows the certification of a nearby system, within a controlled neighborhood of the original one, which attains a single singular point.

The numerical treatment of singular zeroes is a difficult task, mainly because of the ill-posedness of the problem. The following strategy can however be adopted. Find a perturbation of the input system such that the root is a deformation of an exact multiple root. Certainly, there is not a single multiple system, if the input data is approximate. But using the knowledge of the dual structure and interval arithmetic, our method aims at providing a controlled deformation that is compatible with the input.

In this way, we identify the multiplicity structure and we are able to setup deflation techniques which restore the quadratic convergence of the Newton system. The certification of the multiple root is also possible on the symbolically perturbed system by applying a fixed point theorem, based e.g. on interval arithmetic [22] or  $\alpha$ -theorems ([7] and references therein).

This approach has already been explored in the past. The first algebraic work on the analysis of singular points may be due to F.S. Macaulay [14], who introduced the terminology of “inverse system”. His so-called *dialytic method* has been exploited in [12, 13, 4] to construct the inverse system of a multiple point.

Another construction of inverse systems is described e.g. in [17], reducing the size of the intermediate linear systems (and exploited in [23]).

In [18], another approach to construct the dual basis at the singular point which is based on an integration strategy, has been proposed.

Regarding deflation techniques, in [20], by applying a triangulation pre-processing step on the Jacobian matrix at the approximate root, minors of the Jacobian matrix are added to the system to reduce the multiplicity.

In [11], a representation of the ideal in a triangular form in a good position and derivations with respect to the leading variables are used to iteratively reduce the multiplicity. This process is applied for p-adic lifting with exact computation.

In [12, 13], instead of triangulating the Jacobian matrix, the number of variables is doubled and new equations are introduced, which are linear in the new variables. They describe the kernel of the Jacobian matrix at the multiple root. The process is iterative, yet for some practical applications, the root may already be deflated with a few iterations.

In [4], the deflation method is applied iteratively until the root becomes regular, doubling each time the number of variables.

In [21], a minimization approach is used to reduce the value of the equations and their derivatives at the approximate root, assuming a basis of the inverse system is known.

In [24], the inverse system is constructed via Macaulay's method; tables of multiplications are deduced and their eigenvalues are used to improve the approximated root. They show that the convergence is quadratic at the multiple root.

Verification of multiple roots of (approximate) polynomial equations is a difficult task. The approach proposed in [22] consists of introducing perturbation parameters and to certifying the multiple root of nearby system by using a fixed point theorem, based on interval arithmetic. It applies only to cases where the Jacobian has corank equal to 1.

The goal of this paper is to review different techniques that can be used to handle efficiently the following tasks:

- (a) Compute a basis for the dual space and of the local quotient ring at a given (approximate) singular point.
- (b) Deflate the system by augmenting it with new equations derived from the dual basis, introducing adequate perturbation terms.
- (c) Certify the singular point and its multiplicity structure for the perturbed system checking the contraction property of Newton iteration (e.g. via interval arithmetic).

These tools can be applied to improve the quality of approximation of a multiple isolated solution of a system of (polynomial) equations, but they can also be used to solve geometrical problems, such as for instance computing the number of real branches at a singular point of an algebraic curve. For more details on these applications, we refer to [15] and references therein.

## 2 Preliminary considerations

In this section we present some definitions together with the main tools that we shall need in the sequel.

We denote by  $R = \mathbb{K}[\mathbf{x}]$ ,  $\mathbf{x} = (x_1, \dots, x_n)$ , a polynomial ring over the field  $\mathbb{K}$  of characteristic zero. The *dual ring*  $R^*$  is the space of linear functionals  $\Lambda : R \rightarrow \mathbb{K}$ . It is commonly identified to the space of formal series  $\mathbb{K}[[\boldsymbol{\partial}]]$  where  $\boldsymbol{\partial} = (\partial_1, \dots, \partial_n)$  are formal variables. Thus we view dual elements as formal series in differential operators at a point  $\boldsymbol{\zeta} \in \mathbb{K}^n$ . To specify that we use the point  $\boldsymbol{\zeta}$ , we also denote these differentials  $\boldsymbol{\partial}_{\boldsymbol{\zeta}}$ . When applying  $\Lambda(\boldsymbol{\partial}_{\boldsymbol{\zeta}}) \in \mathbb{K}[[\boldsymbol{\partial}_{\boldsymbol{\zeta}}]]$  to a polynomial  $g(\mathbf{x}) \in R$  we will denote by  $\Lambda^{\boldsymbol{\zeta}}[g] = \Lambda^{\boldsymbol{\zeta}}g = \Lambda(\boldsymbol{\partial}_{\boldsymbol{\zeta}})[g(\mathbf{x})]$  the operation

$$\Lambda^{\boldsymbol{\zeta}}[g] = \sum_{\boldsymbol{\alpha} \in \mathbb{N}^n} \frac{\lambda_{\boldsymbol{\alpha}}}{\alpha_1! \cdots \alpha_n!} \cdot \frac{d^{|\boldsymbol{\alpha}|} g}{dx_1^{\alpha_1} \cdots dx_n^{\alpha_n}}(\boldsymbol{\zeta}), \quad (1)$$

for  $\Lambda(\boldsymbol{\partial}_{\boldsymbol{\zeta}}) = \sum \lambda_{\boldsymbol{\alpha}} \frac{1}{\boldsymbol{\alpha}!} \boldsymbol{\partial}_{\boldsymbol{\zeta}}^{\boldsymbol{\alpha}} \in \mathbb{K}[[\boldsymbol{\partial}_{\boldsymbol{\zeta}}]]$ . Extending this definition to an ordered set  $\mathcal{D} = (\Lambda_1, \dots, \Lambda_{\mu}) \in \mathbb{K}[[\boldsymbol{\partial}]]^{\mu}$ , we shall denote  $\mathcal{D}^{\boldsymbol{\zeta}}[g] = (\Lambda_1^{\boldsymbol{\zeta}}g, \dots, \Lambda_{\mu}^{\boldsymbol{\zeta}}g)$ . In some cases, it is convenient to use normalized differentials instead of  $\boldsymbol{\partial}$ : for any  $\boldsymbol{\alpha} \in \mathbb{N}^n$ , we denote  $\mathbf{d}_{\boldsymbol{\zeta}}^{\boldsymbol{\alpha}} = \frac{1}{\boldsymbol{\alpha}!} \boldsymbol{\partial}_{\boldsymbol{\zeta}}^{\boldsymbol{\alpha}}$ . In particular, with the use of this notation we recover the nice property that, if  $\boldsymbol{\zeta} = \mathbf{0}$ , we have  $\mathbf{d}_{\mathbf{0}}^{\boldsymbol{\alpha}} \mathbf{x}^{\boldsymbol{\beta}} = 1$  if  $\boldsymbol{\alpha} = \boldsymbol{\beta}$  and 0 otherwise.

More generally,  $(\mathbf{d}_{\boldsymbol{\zeta}}^{\boldsymbol{\alpha}})_{\boldsymbol{\alpha} \in \mathbb{N}^n}$  is the dual basis of  $((\mathbf{x} - \boldsymbol{\zeta})^{\boldsymbol{\alpha}})_{\boldsymbol{\alpha} \in \mathbb{N}^n}$ , i.e., a non-zero root implies a linear transformation of the variables, so that the root is translated to  $(0, 0)$ .

*Example 1.* Consider the integral of a polynomial function  $g \in R$  over the unit hypercube. Since this is a linear map, it may be expressed in terms of differentials, i.e.:

$$g \mapsto \int_{[0,1]^n} g(\mathbf{x}) dx_1 \cdots dx_n = \sum_{\boldsymbol{\alpha} \in \text{sup}(g)} c_{\boldsymbol{\alpha}} \mathbf{d}^{\boldsymbol{\alpha}}[g],$$

where  $\mathbf{d}^{\boldsymbol{\alpha}}[g] = \frac{1}{\boldsymbol{\alpha}!} \frac{\partial^{|\boldsymbol{\alpha}|} g}{\partial \mathbf{x}^{\boldsymbol{\alpha}}}(\mathbf{0})$  and  $\text{sup}(g)$  stands for the support of  $g$ . Indeed, it can be verified using simple calculations that the (unique) coefficients are given by  $c_{\boldsymbol{\alpha}} = \prod_{i=1}^n \frac{1}{\alpha_i + 1}$ .

For  $\Lambda \in R^*$  and  $p \in R$ , let us define the operation  $p \cdot \Lambda : q \mapsto \Lambda(pq)$ . We check that

$$(x_i - \zeta_i) \cdot \boldsymbol{\partial}_{\boldsymbol{\zeta}}^{\boldsymbol{\alpha}} = \frac{d}{d\partial_{i,\zeta}}(\boldsymbol{\partial}_{\boldsymbol{\zeta}}^{\boldsymbol{\alpha}}), \quad (2)$$

and  $R^*$  obtains the structure of an  $R$ -module. This property shall be useful in the sequel.

## 2.1 Isolated points and differentials

Let  $\mathcal{I} = \langle f_1, \dots, f_s \rangle$  be an ideal of  $R$  and  $\zeta \in \mathbb{K}^n$  a root of the polynomial system  $\mathbf{f} = (f_1, \dots, f_s)$ . We call  $\zeta$  an isolated zero of  $V(\mathcal{I})$  if, in a primary decomposition of  $\mathcal{I}$ , the radical of one of the primary components is the maximal ideal  $m_\zeta = \langle x_1 - \zeta_1, \dots, x_n - \zeta_n \rangle$  defining  $\zeta$  and no other primary component is contained in  $m_\zeta$ .

Suppose that  $\zeta$  is an isolated root of  $\mathbf{f}$ , then a minimal primary decomposition of

$$\mathcal{I} = \bigcap_{\mathcal{Q} \text{ prim. } \supset \mathcal{I}} \mathcal{Q}$$

contains a primary component  $\mathcal{Q}_\zeta$  such that  $\sqrt{\mathcal{Q}_\zeta} = m_\zeta$  and  $\sqrt{\mathcal{Q}'} \not\subset m_\zeta$  for the other primary components  $\mathcal{Q}'$  associated to  $\mathcal{I}$  [2].

As  $\sqrt{\mathcal{Q}_\zeta} = m_\zeta$ , it follows that  $R/\mathcal{Q}_\zeta$  is a finite dimensional vector space. The multiplicity  $\mu_\zeta$  of  $\zeta$  is defined as the dimension of  $R/\mathcal{Q}_\zeta$ . A point of multiplicity one is called *regular point*, or *simple root*, otherwise we say that  $\zeta$  is a singular isolated point, or multiple root of  $\mathbf{f}$ . In the latter case we have  $J_{\mathbf{f}}(\zeta) = 0$ .

*Example 2.* Consider the ideal  $\mathcal{I} = \langle x_1 - x_2 + x_1^2, x_1 - x_2 + x_2^2 \rangle$ , and the root  $\zeta = (0, 0)$ . Then a minimal primary decomposition of  $\mathcal{I}$  is

$$\mathcal{I} = \langle x_2^3, x_1 - x_2 + x_2^2 \rangle \cap \langle -2 + x_2, 2 + x_1 \rangle.$$

Among the two factors we find the maximal ideal of  $\zeta$  given by the radical ideal  $\sqrt{\langle x_2^3, x_1 - x_2 + x_2^2 \rangle} = \langle x_1, x_2 \rangle$ .

We can now define the dual space of an ideal.

**Definition 1** *The dual space of  $\mathcal{I}$  is the subspace of elements of  $\mathbb{K}[[\partial_\zeta]]$  (formal series of the variables  $\partial_\zeta$ ),  $\zeta \in V(\mathcal{I})$ , that vanish on all the elements of  $\mathcal{I}$ . It is also called the orthogonal of  $\mathcal{I}$  and is denoted as  $\mathcal{I}^\perp$ .*

The dual space is known to be isomorphic to the quotient  $R/\mathcal{I}$ . Consider now the *orthogonal* of  $\mathcal{Q}_\zeta$ , i.e. the subspace  $\mathcal{D}_\zeta$  of elements of  $R^*$  that vanish on members of  $\mathcal{Q}_\zeta$ , namely

$$\mathcal{Q}_\zeta^\perp = \mathcal{D}_\zeta = \{\Lambda \in R^* : \Lambda^\zeta[p] = 0, \forall p \in \mathcal{Q}_\zeta\}.$$

The following is an essential property that allows extraction of the local structure  $\mathcal{D}_\zeta$  directly from the “global” ideal  $\mathcal{I} = \langle \mathbf{f} \rangle$ , notably by matrix methods that will be outlined in Section 3.

**Proposition 1** ([18, Th. 8]). *For any isolated point  $\zeta \in \mathbb{K}$  of  $\mathbf{f}$ , we have  $\mathcal{I}^\perp \cap \mathbb{K}[\partial_\zeta] = \mathcal{D}_\zeta$ .*

In other words, we can identify  $\mathcal{D}_\zeta = \mathcal{Q}_\zeta^\perp$  with the space of polynomial differential operators that vanish at  $\zeta$  on every element of  $\mathcal{I}$ . Also note that  $\mathcal{D}_\zeta^\perp = \mathcal{Q}_\zeta$ .

The space  $\mathcal{D}_\zeta$  has dimension  $\mu_\zeta$ , the multiplicity at  $\zeta$ . As the variables  $(x_i - \zeta_i)$  act on  $R^*$  as derivations (see (2)),  $\mathcal{D}_\zeta$  is a space of differential polynomials in  $\mathcal{D}_\zeta$ , which is stable under derivation. This property will be used explicitly in constructing  $\mathcal{D}_\zeta$  (Sect. 3).

**Definition 2** *The nilindex of  $\mathcal{Q}_\zeta$  is the maximal integer  $N \in \mathbb{N}$  such that  $m_\zeta^N \notin \mathcal{Q}_\zeta$ .*

It is directly seen that the maximal order of elements in  $\mathcal{D}_\zeta$  is equal to  $N$ , also known as the *depth* of the space.

## 2.2 Quotient ring and dual structure

In this section we explore the relation between the dual ring and the quotient  $R/\mathcal{Q}_\zeta$ , where  $\mathcal{Q}_\zeta$  is the primary component of the isolated point  $\zeta$ . We show how to extract a basis of this quotient ring from the support of the elements of  $\mathcal{D}_\zeta$  and how  $\mathcal{D}_\zeta$  can be used to reduce any polynomial modulo  $\mathcal{Q}_\zeta$ .

It is convenient in terms of notation to make the assumption  $\zeta = \mathbf{0}$ . This poses no constraint, since it implies only a linear change of coordinates.

Let  $\text{supp } \mathcal{D}_0$  be the set of exponents of monomials appearing in  $\mathcal{D}_0$ , with a non-zero coefficient. These are of degree at most  $N$ , the nilindex of  $\mathcal{Q}_0$ . Since

$$(\forall \Lambda \in \mathcal{D}_0, \Lambda^0[p] = 0) \text{ iff } p \in \mathcal{D}_0^\perp = \mathcal{Q}_0,$$

we derive that  $\text{supp } \mathcal{D}_0 = \{\gamma : \mathbf{x}^\gamma \notin \mathcal{Q}_0\}$ . In particular, we can find a basis of  $R/\mathcal{Q}_0$  between the monomials  $\{\mathbf{x}^\gamma : \gamma \in \text{supp } \mathcal{D}_0\}$ . This is a finite set of monomials, since their degree is bounded by the nilindex of  $\mathcal{Q}_0$ . Now let  $\mathbf{x}^{\gamma_j}, j = 1, \dots, s$  be an enumeration of these monomials. It is clear that these are finitely many, since  $\mathcal{Q}_0$  is zero-dimensional. Given a monomial basis  $\mathcal{B} = (\mathbf{x}^{\beta_i})_{i=1, \dots, \mu}$  of  $R/\mathcal{Q}_0$  and, for all monomials  $\mathbf{x}^{\gamma_j} \notin \mathcal{Q}_0$ , the expression (normal form)

$$\mathbf{x}^{\gamma_j} = \sum_{i=1}^{\mu} \lambda_{ij} \mathbf{x}^{\beta_i} \pmod{\mathcal{Q}_0} \quad (3)$$

of  $\mathbf{x}^{\gamma_j}$  in the basis  $\mathcal{B}$ , then the dual elements [18, Prop. 13]

$$A_i(\mathbf{d}) = \mathbf{d}^{\beta_i} + \sum_{j=1}^{s-\mu} \lambda_{ij} \mathbf{d}^{\gamma_j}, \quad (4)$$

for  $i = 1, \dots, \mu$  form a basis of  $\mathcal{D}_\zeta$ . We give a proof of this fact in the following lemma.

**Lemma 1.** *The set of elements  $\mathcal{D} = (\Lambda_i)_{i=1,\dots,\mu}$  defined in (4) is a basis of  $\mathcal{D}_\zeta$  and the normal form of any  $g(\mathbf{x}) \in R$  with respect to the monomial basis  $\mathcal{B} = (\mathbf{x}^{\beta_i})_{i=1,\dots,\mu}$  is*

$$\text{NF}(g) = \sum_{i=1}^{\mu} \Lambda_i^\zeta[g] \mathbf{x}^{\beta_i}. \quad (5)$$

*Proof.* First note that the elements of  $\mathcal{D}$  are linearly independent, since  $\mathbf{d}^{\beta_i}$  appears only in  $\Lambda_i(\mathbf{d})$ . Now, by construction,

$$\sum_{i=1}^{\mu} \Lambda_i^\zeta[\mathbf{x}^\alpha] \mathbf{x}^{\beta_i} = \text{NF}(\mathbf{x}^\alpha),$$

for all  $\mathbf{x}^\alpha \notin \mathcal{Q}_\zeta$ , e.g.  $\text{NF}(\mathbf{x}^{\beta_i}) = \mathbf{x}^{\beta_i}$ . Also, for  $\mathbf{x}^\alpha \in \mathcal{Q}_\zeta$ ,  $\forall i$ ,  $\Lambda_i^\zeta(\mathbf{x}^\alpha) = 0$ , since  $\alpha \notin \text{supp } \mathcal{D}$ . Thus the elements of  $\mathcal{D}$  compute  $\text{NF}(\cdot)$  on all monomials of  $R$ , and (5) follows by linearity. We deduce that  $\mathcal{D}$  generates the dual, as in Definition 1.  $\square$

It becomes clear that with the knowledge of the dual basis at  $\zeta$ , we are able to compute any  $g \in R$  modulo  $\mathcal{Q}_\zeta$  by applying the basis elements to the monomials of  $g$  (formal derivation plus evaluation at  $\zeta$ ). This lemma also shows an isomorphism between the dual  $\mathcal{D}_\zeta$  and the quotient ring  $R/\mathcal{Q}_\zeta$ , since it implies a one-to-one mapping between the primal and dual basis.

*Example 3.* Consider  $f(x, y) = x^4 + 2x^2y^2 + y^4 + 3x^2y - y^3$  and  $g(x, y) = 18xy^2 - 6x^3$ . The common zero  $\zeta = (0, 0)$  yields the local dual space

$$\mathcal{D} = (1, d_x, d_y, d_x^2, d_x d_y, d_y^2, d_x^3 + \frac{1}{3} d_x d_y^2, d_x^2 d_y + 3 d_y^3, d_x^4 + \frac{1}{3} d_x^2 d_y^2 + d_y^4 + \frac{8}{3} d_y^3),$$

therefore  $\zeta$  is a singular zero with multiplicity  $m = 9$ .

The primal counterpart is  $\mathcal{B} = (1, x, y, x^2, xy, y^2, x^3, x^2y, x^4, xy^2, y^3, x^2y^2, y^4)$ . The relation between  $\mathcal{D}$  and  $\mathcal{B}$  is revealed in the following construction:

$$\begin{array}{c} 1 \quad x \quad y \quad x^2 \quad xy \quad y^2 \quad x^3 \quad x^2y \quad x^4 \quad xy^2 \quad y^3 \quad x^2y^2 \quad y^4 \\ \begin{array}{l} 1 \\ d_x \\ d_y \\ d_x^2 \\ d_x d_y \\ d_y^2 \\ d_x^3 \\ d_x^2 d_y \\ d_x^4 \end{array} \end{array} \begin{bmatrix} 1 & & & & & & & & & 0 & 0 & 0 & 0 \\ & 1 & & & & & & & & 0 & 0 & 0 & 0 \\ & & 1 & & & & & & & 0 & 0 & 0 & 0 \\ & & & 1 & & & & & & 0 & 0 & 0 & 0 \\ & & & & 1 & & & & & 0 & 0 & 0 & 0 \\ & & & & & 1 & & & & 0 & 0 & 0 & 0 \\ & & & & & & 1 & & & 1/3 & 0 & 0 & 0 \\ & & & & & & & 1 & & 0 & 3 & 0 & 0 \\ & & & & & & & & 1 & 0 & 8/3 & 1/3 & 1 \end{bmatrix}. \quad (6)$$

The dual monomial of every row couples with a primal monomial in the corresponding column. From the rows we read the coefficients basis elements in  $\mathcal{D}$ .



The leftmost  $9 \times 9$  block is the identity matrix, implying the duality between  $\mathcal{D}$  and  $\mathcal{B}$ . Then in the last four columns there are some extra monomials; these do not belong to the basis  $\mathcal{B}$ , yet appear with a non-zero coefficient in  $\mathcal{D}$ . These monomials are not in  $\mathcal{Q}_\zeta$ , but they can be reduced modulo  $\mathcal{B}$ : the last four columns yield the normal form of these monomials with respect to  $\mathcal{B}$ . For example, using column 11 we find  $y^3 = 3x^2y + \frac{8}{3}x^4 \pmod{\mathcal{Q}_\zeta}$ .

Using the normal form formula (5), we can derive the table of multiplication by  $x$  and  $y$  in the quotient algebra represented by  $\mathcal{B}$ . To do this, it suffices to compute  $\text{NF}(yx^{\alpha_i}y^{\beta_i})$  and  $\text{NF}(xx^{\alpha_i}y^{\beta_i})$ , for all monomials  $x^{\alpha_i}y^{\beta_i} \in \mathcal{B}$ . This computation can be done by looking up the normal form of each monomial from the rows of matrix (6). The coefficients of these normal forms fill the  $i$ -th rows of the matrices

$$M_x = \begin{bmatrix} 0 & 1 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 1 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 1 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 1 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 1 & 1/3 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 1/3 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 1 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \end{bmatrix} \quad \text{and} \quad M_y = \begin{bmatrix} 0 & 0 & 1 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 1 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 1 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 1/3 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 3/8 & 3/8 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 1/3 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \end{bmatrix}.$$

Computing the normal form of the border monomials of  $\mathcal{B}$  via (5) also yields the border basis relations and the operators of multiplication in the quotient  $R/\mathcal{Q}_\zeta$  (see e.g. [5] for more properties).

If a graded monomial ordering is fixed and  $\mathcal{B} = (\mathbf{x}^{\beta_i})_{i=1,\dots,\mu}$  is the corresponding monomial basis of  $R/\mathcal{Q}_0$ , then  $\mathbf{d}^{\beta_i}$  is the leading term of (4) with respect to the reversed ordering (that is, we reverse the outcome of the comparison of two monomials, keeping equality unchanged) [13, Th. 3.1].

Conversely, if we are given a basis  $\mathcal{D}$  of  $\mathcal{D}_\zeta$  whose coefficient matrix in the dual monomials basis  $(\mathbf{d}^\alpha)_{\alpha \notin \mathcal{Q}_\zeta}$  is  $D \in \mathbb{K}^{\mu \times s}$ , we can compute a basis of  $R/\mathcal{Q}_\zeta$  by choosing  $\mu$  independent columns of  $D$ , say those indexed by  $\mathbf{d}^{\beta_i}$ ,  $i = 1, \dots, \mu$ . If  $G \in \mathbb{K}^{\mu \times \mu}$  is the (invertible) matrix formed by these columns, then  $D' := G^{-1}D$ , is

$$D' = \begin{matrix} & \beta_1 & \cdots & \beta_\mu & \gamma_1 & \cdots & \gamma_{s-\mu} \\ \begin{matrix} A'_1 \\ \vdots \\ A'_\mu \end{matrix} & \begin{bmatrix} 1 & & 0 & \lambda_{1,1} & \cdots & \lambda_{1,s-\mu} \\ & \ddots & & \vdots & & \vdots \\ 0 & & 1 & \lambda_{\mu,1} & \cdots & \lambda_{\mu,s-\mu} \end{bmatrix} \end{matrix}, \quad (7)$$

i.e. a basis of the form (4). Note that an arbitrary basis of  $\mathcal{D}$  does not have the above diagonal form, nor does it directly provide a basis for  $R/Q_\zeta$ . However, a basis of this form has the desired property

$$A_i[\mathbf{x}^{\beta_j}] = \begin{cases} 1, & \text{if } i = j \\ 0, & \text{if } i \neq j \end{cases},$$

for all  $i = 1, \dots, \mu$ .

For  $t \in \mathbb{N}$ ,  $\mathcal{D}_t$  denotes the vector space of polynomials of  $\mathcal{D}$  of degree  $\leq t$ . The Hilbert function  $h : \mathbb{N} \rightarrow \mathbb{N}$  is defined by  $h(t) = \dim(\mathcal{D}_t)$ ,  $t \geq 0$ , hence  $h(0) = 1$  and  $h(t) = \dim \mathcal{D}$  for  $t \geq N$ . The integer  $h(1) - 1 = \text{corank } J_{\mathbf{f}}$  is known as the *breadth* of  $\mathcal{D}$ .

### 3 Computing local ring structure

The computation of a local basis, given a system and a point, is done essentially by matrix-kernel computations, and consequently it can be carried out numerically, even when the point or even the system is inexact. Throughout the section we suppose  $\mathbf{f} \in R^m$  and  $\zeta \in \mathbb{K}^n$  with  $\mathbf{f}(\zeta) = 0$ .

Several matrix constructions have been proposed that use different conditions to identify the dual space as a null-space. They are based on the *stability property* of the dual basis:

$$\forall \Lambda \in \mathcal{D}_t, \quad \frac{d}{d\partial_i} \Lambda \in \mathcal{D}_{t-1}, \quad i = 1, \dots, n. \quad (8)$$

We list existing algorithms that compute dual-space bases:

- As pointed out in (2), an equivalent form of (8) is

$$\forall \Lambda \in \mathcal{D}_t, \quad A[g_i f_i] = 0, \quad \forall g_i \in R \iff A[\mathbf{x}^\beta \cdot f_i] = 0, \quad \forall \beta \in \mathbb{N}^n \quad (9)$$

Macaulay's method [14] uses this equivalent characterization to derive the algorithm that is outlined in Sect. 3.1.

- In [17] they exploit (8) by forming the matrix  $D_i$  of the map

$$\frac{d}{d\partial_i} : \mathbb{K}[\partial]_t \rightarrow \mathbb{K}[\partial]_{t-1}$$

for all  $i = 1, \dots, n$  and some triangular decomposition of the differential polynomials in terms of differential variables. This approach was used in [23] to reduce the row dimension of Macaulay's matrix, but not the column dimension.

- The closedness subspace method of Zeng [25], uses the same condition to identify a superset of  $\text{supp } \mathcal{D}_{t+1}$  when a basis of  $\mathcal{D}_t$  is computed, and thus reduces the column dimension of the matrix.

- The *integration method* in [18] “integrates” elements of a basis of  $\mathcal{D}_t$ , and obtains *a priori* knowledge of the form of elements in degree  $t + 1$  (Sect. 3.2).

All methods are incremental, in the sense that they start by setting  $\mathcal{D}_0 = (\mathbf{1})$  and continue by computing  $\mathcal{D}_i$ ,  $i = 1, \dots, N, N + 1$ . When  $\#\mathcal{D}_N = \#\mathcal{D}_{N+1}$  then  $\mathcal{D}_N$  is a basis of  $\mathcal{D}$ , and  $N$  is the nilindex of  $\mathcal{Q}$ .

We shall review two of these approaches to compute a basis for  $\mathcal{D}$ , and then describe an improvement, that allows simultaneous computation of a monomial basis of the quotient ring, while avoiding redundant computations.

### 3.1 Macaulay’s dialytic matrices

This matrix construction is presented in [14, Ch. 4], a modern introduction is contained in [4], together with an implementation of the method in ApaTools<sup>1</sup>.

The idea behind the algorithm is the following: An element of  $\mathcal{D}$  is of the form

$$\Lambda(d) = \sum_{|\alpha| \leq N} \lambda_\alpha d^\alpha$$

under the condition:  $\Lambda^0$  evaluates to 0 at any  $g \in \langle \mathbf{f} \rangle$ , that is,

$$\Lambda^0(g) = \Lambda^0\left(\sum g_i f_i\right) = 0 \iff \Lambda^0(x^\beta f_i) = 0 \quad ,$$

for all monomials  $x^\beta$ ,  $\beta \in \mathbb{N}$ .

If we apply this condition recursively for  $|\alpha| \leq N$ , we get a vector of coefficients  $(\lambda_\alpha)_{|\alpha| \leq N}$  in the (right) kernel of the matrix with rows indexed by constraints  $\Lambda^0[x^\beta f_i] = 0$ ,  $|\beta| \leq N - 1$ . A basis of  $\mathcal{D}_N$  is given by the kernel of this matrix in depth  $N$ . The method consists in computing the kernel of these matrices for  $N = 1, 2, \dots$ ; when  $N$  reaches the nilindex of  $\mathcal{I}$ , For some value of  $N$ , this kernel stabilizes and the generating vectors form a basis of  $\mathcal{D}$ .

Note that the only requirement is to be able to perform derivation of the input equations and evaluation at  $\zeta = \mathbf{0}$ .

*Example 4.* Let  $f_1 = x_1 - x_2 + x_1^2$ ,  $f_2 = x_1 - x_2 + x_2^2$ . We also refer the reader to [4, Ex. 2] for a detailed demonstration of this instance. The matrices in order 1 and 2 are:

$$\begin{array}{c} \begin{array}{ccc} 1 & d_1 & d_2 \\ f_1 & \begin{bmatrix} 0 & 1 & -1 \end{bmatrix} \\ f_2 & \begin{bmatrix} 0 & 1 & -1 \end{bmatrix} \end{array} , \quad \begin{array}{c} \begin{array}{cccccc} 1 & d_1 & d_2 & d_1^2 & d_1 d_2 & d_2^2 \\ f_1 & 0 & 1 & -1 & 1 & 0 & 0 \\ f_2 & 0 & 1 & -1 & 0 & 0 & 1 \\ x_1 f_1 & 0 & 0 & 0 & 1 & -1 & 0 \\ x_1 f_2 & 0 & 0 & 0 & 1 & -1 & 0 \\ x_2 f_1 & 0 & 0 & 0 & 0 & 1 & -1 \\ x_2 f_2 & 0 & 0 & 0 & 0 & 1 & -1 \end{array} \end{array} \end{array} .$$

<sup>1</sup> <http://www.neiu.edu/~zzeng/apatools.htm>

The kernel of the left matrix gives  $\mathcal{D}_1 = (1, d_1 + d_2)$ . Expanding up to order two, we get the matrix on the right, and  $\mathcal{D}_2 = (1, d_1 + d_2, -d_1 + d_1^2 + d_1 d_2 + d_2^2)$ . If we expand up to depth 3 we get the same null-space, thus  $\mathcal{D} = \mathcal{D}_2$ .

### 3.2 Integration method

This method is presented in [18]. It is an evolution of Macaulay's method, since the matrices are not indexed by all differentials, but just by elements based on knowledge of the previous step. This performs a computation adapted to the given input and results in smaller matrices.

For  $\Lambda \in \mathbb{K}[\partial]$ , we denote by  $\int_k \Lambda$  the element  $\Phi \in \mathbb{K}[\partial]$  with the property  $\frac{d}{d\partial_k} \Phi(\partial) = \Lambda(\partial)$  and with no constant term with respect to  $\partial_k$ .

**Theorem 3** ([18, Th. 15]). *Let  $\{\Lambda_1, \Lambda_2, \dots, \Lambda_s\}$  be a basis of  $\mathcal{D}_{t-1}$ , that is, the subspace of  $\mathcal{D}$  of elements of order at most  $t-1$ . An element  $\Lambda \in \mathbb{K}[\partial]$  with no constant term lies in  $\mathcal{D}_t$  iff it is of the form:*

$$\Lambda(\partial) = \sum_{i=1}^s \sum_{k=1}^n \lambda_{ik} \int_k \Lambda_i(\partial_1, \dots, \partial_k, 0, \dots, 0), \quad (10)$$

for  $\lambda_{ik} \in \mathbb{K}$ , and the following two conditions hold:

$$(i) \sum_{i=1}^s \lambda_{ik} \frac{d}{d\partial_l} \Lambda_i(\partial) - \sum_{i=1}^s \lambda_{il} \frac{d}{d\partial_k} \Lambda_i(\partial) = 0, \text{ for all } 1 \leq k < l \leq n.$$

$$(ii) \Lambda^\zeta[f_k] = 0, \text{ for } k = 1, \dots, m.$$

Condition (i) is equivalent to  $\frac{d}{d\partial_k} \Lambda \in \mathcal{D}_{t-1}$ , for all  $k$ . Thus the two conditions express exactly the fact that  $\mathcal{D}$  must be stable under derivation and its members must vanish on  $\langle f \rangle$ .

This gives the following algorithm to compute the dual basis: Start with  $\mathcal{D}_0 = \langle 1 \rangle$ . Given a basis of  $\mathcal{D}_{t-1}$  we generate the  $ns$  candidate elements  $\int_k \Lambda_{i-1}(\partial_1, \dots, \partial_k, 0, \dots, 0)$ . Conditions (i) and (ii) give a linear system with unknowns  $\lambda_{ik}$ . The columns of the corresponding matrix are indexed by the candidate elements. Then, the kernel of this matrix gives a basis of  $\mathcal{D}_t$ , which we use to generate new candidate elements. If for some  $t$  we compute a kernel of the same dimension as  $\mathcal{D}_{t-1}$ , then we have a basis of  $\mathcal{D}$ .

*Example 5.* Consider the instance of Ex. 4,  $f_1 = x_1 - x_2 + x_1^2$ ,  $f_2 = x_1 - x_2 + x_2^2$ . We have  $f_1(\zeta) = f_2(\zeta) = 0$ , thus we set  $\mathcal{D}_0 = \{1\}$ . Equation (10) gives  $\Lambda = \lambda_1 d_1 + \lambda_2 d_2$ . Condition (i) induces no constraints and (ii) yields the system

$$\begin{bmatrix} 1 & -1 \\ 1 & -1 \end{bmatrix} \begin{bmatrix} \lambda_1 \\ \lambda_2 \end{bmatrix} = 0 \quad (11)$$

where the columns are indexed by  $d_1, d_2$ . We get  $\lambda_1 = \lambda_2 = 1$  from the kernel of this matrix, thus  $\mathcal{D}_1 = \{1, d_1 + d_2\}$ .

For the second step, we compute the elements of  $\mathcal{D}_2$ , that must be of the form

$$A = \lambda_1 d_1 + \lambda_2 d_2 + \lambda_3 d_1^2 + \lambda_4 (d_1 d_2 + d_2^2).$$

Condition (i) yields  $\lambda_3 - \lambda_4 = 0$ , and together with (ii) we form the system

$$\begin{bmatrix} 0 & 0 & 1 & -1 \\ 1 & -1 & 1 & 0 \\ 1 & -1 & 0 & 1 \end{bmatrix} \begin{bmatrix} \lambda_1 \\ \vdots \\ \lambda_4 \end{bmatrix} = 0, \quad (12)$$

with columns indexed by  $d_1, d_2, d_1^2, d_1 d_2 + d_2^2$ . We get two vectors in the kernel, the first yielding again  $d_1 + d_2$  and a second one for  $\lambda_1 = -1, \lambda_2 = 0, \lambda_3 = \lambda_4 = 1$ , so we deduce that  $-d_1 + d_1^2 + d_1 d_2 + d_2^2$  is a new element of  $\mathcal{D}_2$ .

In the third step we have

$$A = \lambda_1 d_1 + \lambda_2 d_2 + \lambda_3 d_1^2 + \lambda_4 (d_1 d_2 + d_2^2) + \lambda_5 (d_1^3 - d_1^2) + \lambda_6 (d_2^3 + d_1 d_2^2 + d_1^2 d_2 - d_1 d_2), \quad (13)$$

condition (i) leads to  $\lambda_3 - \lambda_4 + (\lambda_5 - \lambda_6)(d_1 + d_2) = 0$ , and together with condition (ii) we arrive at

$$\begin{bmatrix} 0 & 0 & 0 & 0 & 1 & -1 \\ 0 & 0 & 1 & -1 & 0 & 0 \\ 1 & -1 & 1 & 0 & -1 & 0 \\ 1 & -1 & 0 & 1 & 0 & 0 \end{bmatrix} \begin{bmatrix} \lambda_1 \\ \vdots \\ \lambda_6 \end{bmatrix} = 0, \quad (14)$$

of size  $4 \times 6$ , having two kernel elements that are already in  $\mathcal{D}_2$ . We derive that  $\mathcal{D} = \langle \mathcal{D}_2 \rangle = \langle \mathcal{D}_3 \rangle$  and the algorithm terminates.

Note that for this example Macaulay's method ends with a matrix of size  $12 \times 10$ , instead of  $4 \times 6$  in this approach.

### 3.3 Computing a primal-dual pair

In this section we provide a process that allows simultaneous computation of a basis pair  $(\mathcal{D}, \mathcal{B})$  of  $\mathcal{D}$  and  $R/\mathcal{Q}$ .

Computing a basis of  $\mathcal{D}$  degree by degree involves duplicated computations. The successive spaces computed are  $\mathcal{D}_1 \subset \dots \subset \mathcal{D}_N = \mathcal{D}_{N+1}$ . It is more efficient to produce only new elements  $A \in \mathcal{D}_t$ , independent in  $\mathcal{D}_t/\mathcal{D}_{t-1}$ , at step  $t$ .

Also, once a dual basis is computed, one has to transform it to the form (4), in order to identify a basis of  $R/\mathcal{Q}$  as well. This transformation can be done *a posteriori*, by finding a sub-matrix of full rank and then performing Gauss-Jordan elimination over this sub-matrix, to reach matrix form (7).

We introduce a condition (iii) extending Th. 3, that addresses these two issues: It allows the computation of a total of  $\mu$  independent elements throughout execution, and returns a “triangular” basis, e.g. a basis of  $R/\mathcal{Q}$  is identified.

**Lemma 2.** Let  $\mathcal{D}_{t-1} = (\Lambda_1, \dots, \Lambda_k)$  be a basis of  $\mathcal{D}_{t-1}$ , whose coefficient matrix is

$$\begin{array}{c} \beta_1 \quad \cdots \quad \beta_k \quad \gamma_1 \quad \cdots \quad \gamma_{s-k} \\ \Lambda_1 \begin{bmatrix} 1 & * & * & * & \cdots & * \\ \vdots & \ddots & * & \vdots & \vdots & \vdots \\ 0 & 0 & 1 & * & \cdots & * \end{bmatrix}, \end{array} \quad (15)$$

yielding the monomial basis  $\mathcal{B}_{t-1} = (\mathbf{x}^{\beta_i})_{i=1, \dots, k}$ . Also, let  $\Lambda \in \mathbb{K}[\mathcal{D}]$  be of the form (10), satisfying (i-ii) of Th. 3.

If we impose the additional condition:

$$(iii) \quad \Lambda^\zeta[\mathbf{x}^{\beta_i}] = 0, \quad 1 \leq i \leq k,$$

then the kernel of the matrix implied by (i-iii) is isomorphic to  $\mathcal{D}_t/\mathcal{D}_{t-1}$ . Consequently, it extends  $\mathcal{D}_{t-1}$  to a basis of  $\mathcal{D}_t$ .

*Proof.* Let  $S$  be the kernel of the matrix implied by (i-iii), and let  $\Lambda \in \mathbb{K}[\mathcal{D}]$  be a non-zero functional in  $S$ . We have  $\Lambda \in \mathcal{D}_t$  and  $\Lambda^\zeta[\mathbf{x}^{\beta_i}] = 0$  for  $i = 1, \dots, k$ .

First we show that  $\Lambda \notin \mathcal{D}_{t-1}$ . If  $\Lambda \in \mathcal{D}_{t-1}$ , then  $\Lambda = \sum_{i=1}^k \lambda_i \Lambda_i$ . Take for  $i_0$  the minimal  $i$  such that  $\lambda_i \neq 0$ . Then  $\Lambda^\zeta[\mathbf{x}^{\beta_{i_0}}] = \lambda_{i_0}$ , which contradicts condition (iii). Therefore,  $S \cap \mathcal{D}_{t-1} = \{0\}$ , and  $S$  can be naturally embedded in  $\mathcal{D}_t/\mathcal{D}_{t-1}$ , i.e.  $\dim S \leq \dim \mathcal{D}_t - \dim \mathcal{D}_{t-1}$ .

It remains to show that  $\dim S$  is exactly  $\dim \mathcal{D}_t - \dim \mathcal{D}_{t-1}$ . This is true, since with condition (iii) we added  $k = \dim \mathcal{D}_{t-1}$  equations, thus we excluded from the initial kernel (equal to  $\mathcal{D}_t$ ) of (i-ii) a subspace of dimension at most  $k = \dim \mathcal{D}_{t-1}$ , so that  $\dim S \geq \dim \mathcal{D}_t - \dim \mathcal{D}_{t-1}$ .

We deduce that  $S \cong \mathcal{D}_t/\mathcal{D}_{t-1}$ , thus a basis of  $S$  extends  $\mathcal{D}_{t-1}$  to a basis of  $\mathcal{D}_t$ .  $\square$

The above condition is easy to realize; it is equivalent to  $\forall i, \mathbf{d}^{\beta_i} \notin \text{supp } \Lambda$ , which implies adding a row (linear constraint) for every  $i$ . If we choose the elements of  $\mathcal{B}$  with a “reversed” total degree ordering (if a monomial compares total-degree “less than” another one, then it compares “bigger than” the same monomial in the reversed order), then in many cases this constraint becomes  $\lambda_{ik} = 0$  for some  $i, k$ . In this case we rather remove the column corresponding to  $\lambda_{ik}$  instead of adding a row. Hence this lemma allows to shrink the kernel (but also the dimension) of the matrix and compute only new dual elements, which are reduced modulo the previous basis. For a detailed size comparison, see 1.

Let us explore our running example, to demonstrate the essence of this improvement.

*Example 6.* We re-run Ex. 5 using Lem. 2. In the initialization step  $\mathcal{D}_0 = (1)$  is already in triangular form with respect to  $\mathcal{B}_0 = \{1\}$ . For the first step, we demand  $\Lambda[1] = 0$ , thus the matrix is the same as (11), yielding  $\mathcal{D}_1 =$

$(1, d_1 + d_2)$ . We extend  $\mathcal{B}_1 = \{1, x_2\}$ , so that  $\mathcal{D}_1$  is triangular with respect to  $\mathcal{B}_1$ .

In the second step we remove from [12](#) the second column, hence we are left with the  $3 \times 3$  system

$$\begin{bmatrix} 0 & 1 & -1 \\ 1 & 1 & 0 \\ 1 & 0 & 1 \end{bmatrix} \begin{bmatrix} \lambda_1 \\ \lambda_3 \\ \lambda_4 \end{bmatrix} = 0,$$

yielding a single solution  $-d_1 + d_1^2 + d_1 d_2 + d_2^2$ . We extend  $\mathcal{B}_1$  by adding the monomial  $x_1$ :  $\mathcal{B}_2 = \{1, x_2, x_1\}$ .

For the final step, we search an element of the form [\(13\)](#) with  $A[x_1] = A[x_2] = 0$ , and together with (i–ii) we get:

$$\begin{bmatrix} 0 & 0 & 1 & -1 \\ 1 & -1 & 0 & 0 \\ 1 & 0 & -1 & 0 \\ 0 & 1 & 0 & 0 \end{bmatrix} \begin{bmatrix} \lambda_3 \\ \vdots \\ \lambda_6 \end{bmatrix} = 0.$$

We find an empty kernel, thus we recover the triangular basis  $\mathcal{D} = \mathcal{D}_2$ , which can be diagonalized to reach the form:

$$\begin{matrix} & 1 & d_2 & d_1 & d_1^2 & d_1 d_2 & d_2^2 \\ \begin{matrix} A_1 \\ A_2 \\ A_3 \end{matrix} & \begin{bmatrix} 1 & 0 & 0 & 0 & 0 & 0 \\ 0 & 1 & 0 & 1 & 1 & 1 \\ 0 & 0 & 1 & -1 & -1 & -1 \end{bmatrix} \end{matrix}.$$

This diagonal basis is dual to the basis  $\mathcal{B} = (1, x_2, x_1)$  of the quotient ring and also provides a normal form algorithm (Lem. [1](#)) with respect to  $\mathcal{B}$ . In the final step we generated a  $4 \times 4$  matrix, of smaller size compared to all previous methods.

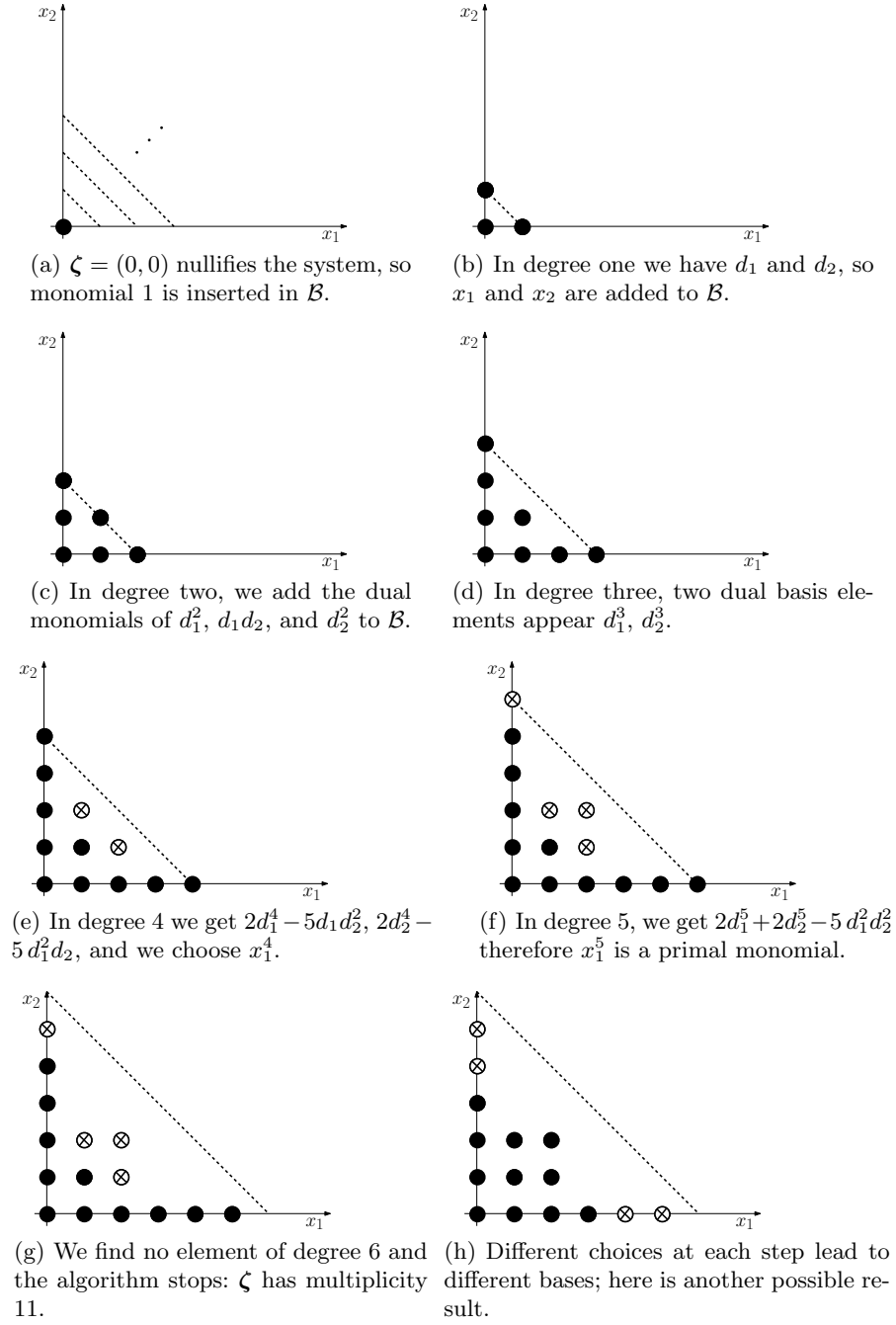
Another example is treated in figure [1](#), with the aid of pictures.

This technique for computing  $\mathcal{B}$  can be applied similarly to other matrix methods, e.g. Macaulay's dialytic method.

If  $h(t) - h(t-1) > 1$ , i.e. there is more than one element in step  $t$ , then the choice of monomials to add to  $\mathcal{B}$  is obtained by extracting a non-zero maximal minor from the coefficient matrix in  $(d^\alpha)$ . In practice, we will look first at the minimum monomials with respect to a fixed term ordering.

## 4 Deflation of a singular point

Deflation techniques allow to transform a system of equations defining a singular solution into a new system where the solution corresponds to a simple point. Usually this is done by adding new variables and new equations so that a simple isolated solution of the extended system projects onto the singular solution of the initial system. We will illustrate different types of deflation.



**Fig. 1.** Discovering a primal-dual basis pair for the root  $\zeta = (0, 0)$  of the bivariate system  $\{2x_1x_2^2 + 5x_1^4, 2x_1^2x_2 + 5x_2^4\}$ .



#### 4.1 The univariate case

In preparation for the multivariate case, we review an approach for treating singularities of univariate polynomials.

Let  $g(x) \in \mathbb{K}[x]$  be a polynomial which attains at  $x = 0$  a root of multiplicity  $\mu > 1$ . The latter is defined as the positive integer  $\mu$  such that  $d^\mu g(0) \neq 0$  whereas  $g(0) = dg(0) = \dots = d^{\mu-1}g(0) = 0$ . Here we denote by  $d^k g(x) = \frac{1}{k!} \frac{d^k}{dx^k} g(x)$  the normalized derivative of order  $k$  with respect to  $x$ .

We see that  $\mathcal{D}_0 = \langle 1, d, \dots, d^{\mu-1} \rangle$  is the maximal space of differentials which is stable under derivation, that vanish when applied to members of  $\mathcal{Q}_0$ , the  $\langle x \rangle$ -primary component of  $\langle g \rangle$  at  $x = 0$ .

*Example 7.* Let  $g(x) = (x - 1)^4$ , and  $\zeta = 1$ . First we check that the space of differentials that vanish on the solution include all linear combinations of  $\mathcal{D}_\zeta = \langle 1, d_\zeta, d_\zeta^2, d_\zeta^3 \rangle$ . For instance, we compute  $d^2[g] = 6(x - 1)^2|_\zeta = 0$ .

Now  $d^4$  is not a member of  $\mathcal{D}_\zeta$  since  $d^4[g] = 1 \neq 0$  does not vanish. Similarly, for all  $i \geq 4$ , there exists a member of the ideal generated by  $g$  which does not evaluate to zero when we apply the differential  $d^i$ , namely  $d^i[x^{i-4}g] = 1$ .

We conclude that the local dual space is exactly  $\mathcal{D}_\zeta$ , and verify that  $\zeta$  4-fold zero of  $f$ .

Consider now the symbolically perturbed equation

$$f_1(x, \varepsilon) = g(x) + \varepsilon_1 + \varepsilon_2 x + \dots + \varepsilon_{\mu-1} x^{\mu-2} \quad (16)$$

and apply every basis element of  $\mathcal{D}_0$  to arrive at the new system

$$\mathbf{f}(x, \varepsilon) = (f_1, df_1, \dots, d^{\mu-1} f_1)$$

in  $\mu - 1$  variables. The  $i$ -th equation has the form

$$f_i = d^{i-1} f_1 = d^{i-1} g + \sum_{k=i}^{\mu-1} \binom{k-1}{i-1} x^{k-i} \varepsilon_k, \quad$$

i.e linear in  $\varepsilon$ , the last one being  $f_\mu = d^{\mu-1} g(x)$ . This system deflates the root, since the determinant of its Jacobian matrix at  $(0, \mathbf{0})$  is

$$\det J_{\mathbf{f}}(0, \mathbf{0}) = \left| \begin{array}{c|cc} \frac{d}{dx} f_1 & 1 & 0 \\ \vdots & \ddots & \\ \frac{d}{dx} f_{\mu-1} & 0 & 1 \\ \hline \frac{d}{dx} f_\mu & 0 & \end{array} \right| = \begin{aligned} &= -df_\mu(0) \\ &= -\mu d^\mu g(0) \neq 0. \end{aligned}$$

Now suppose that  $\zeta^*$  is an approximate zero, close to  $x = \zeta$ . We can still compute  $\mathcal{D}_\zeta$  by evaluating  $g(x)$  and the derivatives up to a threshold relative to the error present in  $\zeta^*$ . Then we can form (16) and use verification

techniques to certify the root. Checking that the Newton operator is contract-  
ing shows the existence and unicity of a multiple root in a neighborhood of  
the input data. We are going to extend this approach, described in [22], to  
multi-dimensional isolated multiple roots.

## 4.2 Deflation using the dialytic approach

Let us consider a system of equations  $\mathbf{f} = (f_1, \dots, f_s)$ ,  $f_k \in \mathbb{R}[\mathbf{x}]$ ,  $s \geq n$ ,  
which has an isolated zero  $\zeta$ .

If the Jacobian matrix

$$J_{\mathbf{f}}(\mathbf{x}) = \begin{pmatrix} \partial_{x_1} f_1(\mathbf{x}) & \cdots & \partial_{x_n} f_1(\mathbf{x}) \\ \vdots & & \vdots \\ \partial_{x_1} f_s(\mathbf{x}) & \cdots & \partial_{x_n} f_s(\mathbf{x}) \end{pmatrix}$$

at the point  $\zeta$  is of (maximal) rank  $n$ , then the root  $\zeta$  is simple. Moreover the  
iteration

$$\mathbf{x}^{(n+1)} = \mathbf{x}^{(n)} + J_{\mathbf{f}}(\mathbf{x}^{(n)})^+ \mathbf{f}(\mathbf{x}^{(n)})$$

converges to  $\zeta$  as soon as the initial point  $\mathbf{x}^{(0)}$  is close enough to  $\zeta$  [9].

When the root is not simple, then the rank of  $J_{\mathbf{f}}(\zeta)$  is  $r_1 < n$  and there  
are  $n - r_1$  linearly independent differentials of order 1 of the form  $u_1 \partial_{x_1} +$   
 $\cdots + u_n \partial_{x_n}$  which satisfy

$$u_1 \partial_{x_1} \mathbf{f}(\zeta) + \cdots + u_n \partial_{x_n} \mathbf{f}(\zeta) = 0, \quad (17)$$

or equivalently

$$J_{\mathbf{f}}(\zeta) \begin{pmatrix} u_1 \\ \vdots \\ u_n \end{pmatrix} = \mathbf{0}.$$

To fix a solution of this system, we can choose  $n - r_1$  random vectors  $\mathbf{r}_j =$   
 $(r_{0,j}, r_{1,j}, \dots, r_{j,n})$  for  $j = 1, \dots, n - r_1$  and consider the equations

$$u_1 r_{1,j} + \cdots + u_n r_{j,n} + r_{j,0} = 0, \quad j = 1, \dots, n - r_1. \quad (18)$$

Extending the initial system of equations  $\mathbf{f}(\mathbf{x})$  with the new equations (17),  
(18), we obtain a new system of equations denoted  $\mathbf{f}_{[1]}(\mathbf{x}_{[1]})$  in the variables  
 $\mathbf{x}_{[1]} = (x_1, \dots, x_n, u_1, \dots, u_n)$ . This system is called a system *deflated* from  $\mathbf{f}$ .

By construction, if  $\zeta$  is an isolated root of  $\mathbf{f}(\mathbf{x}) = 0$  and  $\text{rank } J_{\mathbf{f}}(\zeta) = r_1$ ,  
there is a unique  $\mathbf{u}^*$  satisfying equations (17), (18). Thus  $\mathbf{x}_{[1]}^* = (\zeta, \mathbf{u}^*)$  is an  
isolated solution of the system  $\mathbf{f}_{[1]}(\mathbf{x}_{[1]}) = 0$ .

If the root  $\mathbf{x}_{[1]}^*$  of this system is simple, then Newton iteration applied on  
 $\mathbf{f}_{[1]}$  will converge quadratically to  $\mathbf{x}_{[1]}^*$  for an initial point  $\mathbf{x}_{[1]}^{(0)}$  in its neighbor-  
hood.

If the root is not simple, the deflation can be applied to the system  $\mathbf{f}_{[1]}(\mathbf{x}_{[1]}) = 0$  and we obtain a new system  $\mathbf{f}_{[2]}(\mathbf{x}_{[2]}) = 0$  in  $4n$  variables.

As shown in the next result, the process can be applied inductively until the root becomes simple:

**Theorem 4** ([12], [3]). *If  $\zeta$  is an isolated root of the system  $\mathbf{f}(\mathbf{x}) = 0$ , there exists a number  $k \in \mathbb{N}$  such that  $\mathbf{f}_{[k]}(\mathbf{x}_{[k]}) = 0$  has a simple root  $\mathbf{x}_{[k]}^*$  whose projection on the first  $n$  coordinates is  $\zeta$ .*

It is proved in [12] (or in [3]) that the number  $k$  of iterations is at most the depth of the multiplicity of  $\mathbf{f}$  at  $\zeta$ , that is the maximum degree of a differential polynomial of the inverse system of  $\mathbf{f}$  at  $\zeta$ .

Notice that the number of variables of the system  $\mathbf{f}_{[k]}$  is  $n \times 2^k$ .

*Example 8.* Consider the system  $f_1(x_1, x_2) = 2x_1x_2^2 + 5x_1^4$ ,  $f_2(x_1, x_2) = 2x_1^2x_2 + 5x_2^4$  and the singular point  $\zeta = (0, 0)$ . Since  $J_{\mathbf{f}}(\zeta) = 0$ , we apply a first deflation step, i.e. we compute the equations  $J_{\mathbf{f}}(\mathbf{x}) \begin{pmatrix} u_1 \\ u_2 \end{pmatrix}$ , and two random linear equations:

$$\begin{aligned} g_1(\mathbf{x}_{[1]}) &= f_1 = 2x_1x_2^2 + 5x_1^4, & g_2(\mathbf{x}_{[1]}) &= f_2 = 2x_1^2x_2 + 5x_2^4 \\ g_3(\mathbf{x}_{[1]}) &= (2x_2^2 + 20x_1^3)u_1 + 4x_1x_2u_2 \\ g_4(\mathbf{x}_{[1]}) &= 4x_1x_2u_1 + (2x_1^2 + 20x_2^3)u_2 \\ g_5(\mathbf{x}_{[1]}) &= 16u_1 + u_2 - 1, & g_6(\mathbf{x}_{[1]}) &= 70u_1 + 77u_2 \end{aligned}$$

The new Jacobian matrix  $J_{\mathbf{g}}(\zeta_{[1]})$  is rank-defect, with  $\zeta_{[1]} = (\mathbf{x}, \mathbf{u}) = \left(0, 0, \frac{11}{166}, \frac{-5}{83}\right)$  is zero, and the multiplicity has dropped from 11 to 4. Therefore we repeat the procedure for this new system using random equations from the kernel of  $J_{\mathbf{g}}(\zeta_{[1]})$  (18):

$$\begin{aligned} h_1(\mathbf{x}_{[2]}) &= 2x_1x_2^2 + 5x_1^4, & h_2(\mathbf{x}_{[2]}) &= 2x_1^2x_2 + 5x_2^4 \\ h_3(\mathbf{x}_{[2]}) &= (2x_2^2 + 20x_1^3)u_1 + 4x_1x_2u_2 \\ h_4(\mathbf{x}_{[2]}) &= 4x_1x_2u_1 + (2x_1^2 + 20x_2^3)u_2 \\ h_5(\mathbf{x}_{[2]}) &= 16u_1 + u_2 - 1, & h_6(\mathbf{x}_{[2]}) &= 70u_1 + 77u_2 \\ h_7(\mathbf{x}_{[2]}) &= (2x_2^2 + 20x_1^3)v_1 + 4x_1x_2v_2, & h_8(\mathbf{x}_{[2]}) &= 4x_1x_2v_1 + (2x_1^2 + 20x_2^3)v_2 \\ h_9(\mathbf{x}_{[2]}) &= (60x_1^2u_1 + 4x_2u_2)v_1 + (4x_2u_1 + 4x_1u_2)v_2 + (2x_2^2 + 20x_1^3)v_3 + 4x_1x_2v_4 \\ h_{10}(\mathbf{x}_{[2]}) &= (4x_2u_1 + 4x_1u_2)v_1 + (4x_1u_1 + 60x_2^2u_2)v_2 + 4x_1x_2v_3 + (2x_1^2 + 20x_2^3)v_4 \\ h_{11}(\mathbf{x}_{[2]}) &= 16v_3 + v_4, & h_{12}(\mathbf{x}_{[2]}) &= 70v_3 + 77v_4 \\ h_{13}(\mathbf{x}_{[2]}) &= 53v_1 + 12v_2 + 19v_3 + 63v_4 - 1, & h_{14}(\mathbf{x}_{[2]}) &= 40v_1 + 90v_2 + 3v_3 + 49v_4 \end{aligned}$$

We obtain a new system which has a regular root at  $\zeta_{[2]} = (\mathbf{x}, \mathbf{u}, \mathbf{v}) = \left(0, 0, \frac{11}{166}, \frac{-5}{83}, \frac{3}{143}, \frac{-4}{429}, 0, 0\right)$ , so deflation is achieved in two steps. We see

that the number of (equations and) variables increased exponentially, from 2 to 8, and the system is no longer square.

### 4.3 Deflation using the inverse system

We consider again a system of equations  $\mathbf{f} = (f_1, \dots, f_s)$ ,  $f_k \in \mathbb{R}[\mathbf{x}]$ , which has an isolated root  $\zeta$  of multiplicity  $\mu$ .

In this section, we will also extend the initial system by introducing new variables so that the extended system has a simple isolated root, which projects onto the multiple point  $\zeta$ . Contrarily to the deflation technique described in Section 4.2, the number of new variables will be directly related with the multiplicity  $\mu$  of the point. Let  $\mathbf{b} = ((\mathbf{x} - \zeta)^{\beta_1}, \dots, (\mathbf{x} - \zeta)^{\beta_\mu})$  be a basis of  $R/\mathcal{Q}_\zeta$  and  $\mathcal{D} = (\Lambda_1, \dots, \Lambda_\mu)$  its dual counterpart, with  $\beta_1 = (0, \dots, 0)$ ,  $\Lambda_1 = \mathbf{1}$ .

We introduce a new set of equations starting from  $\mathbf{f}$ , as follows: add to every  $f_k$  the polynomial  $g_k = f_k + p_k$ ,  $p_k = \sum_{i=1}^{\mu} \varepsilon_{i,k} (\mathbf{x} - \zeta)^{\beta_i}$  where  $\varepsilon_k = (\varepsilon_{k,1}, \dots, \varepsilon_{k,\mu})$  is a new vector of  $\mu$  variables.

Consider the system

$$\mathcal{D}\mathbf{g}(\mathbf{x}, \varepsilon) = (\Lambda_1(\partial_{\mathbf{x}})[g], \dots, \Lambda_\mu(\partial_{\mathbf{x}})[g]).$$

where  $\Lambda^{\mathbf{x}}[g_k] = \Lambda_i(\mathbf{d}_{\mathbf{x}})[g_k]$  is defined as in (1) with  $\zeta$  replaced by  $\mathbf{x}$ , i.e. we differentiate  $g_k$  but we do not evaluate at  $\zeta$ . This is a system of  $\mu s$  equations, which we shall index  $\mathcal{D}\mathbf{g}(\mathbf{x}, \varepsilon) = (g_{1,1}, \dots, g_{\mu,s})$ . We have

$$g_{ik}(\mathbf{x}, \varepsilon) = \Lambda_i^{\mathbf{x}}[f_k + p_k] = \Lambda_i^{\mathbf{x}}[f_k] + \Lambda_i^{\mathbf{x}}[p_k] = \Lambda_i^{\mathbf{x}}[f_k] + p_{i,k}(\mathbf{x}, \varepsilon).$$

Notice that  $p_{i,k}(\zeta, \varepsilon) = \Lambda_i^{\zeta}[p_k] = \varepsilon_{i,k}$  because  $\mathcal{D} = (\Lambda_1, \dots, \Lambda_\mu)$  is dual to  $\mathbf{b}$ .

As the first basis element of  $\mathcal{D}$  is  $\mathbf{1}$  (the evaluation at the root), the first  $s$  equations are  $\mathbf{g}(\mathbf{x}, \varepsilon) = 0$ .

Note that this system is under-determined, since the number of variables is  $\mu s + n$  and the number of equations is  $\mu s$ . We shall provide a systematic way to choose  $n$  variables and purge them (or better, set them equal to zero).

By linearity of the Jacobian matrix we have

$$\begin{aligned} J_{\mathcal{D}\mathbf{g}}(\mathbf{x}, \varepsilon) &= J_{\mathcal{D}\mathbf{f}}(\mathbf{x}, \varepsilon) + J_{\mathcal{D}\mathbf{p}}(\mathbf{x}, \varepsilon) \\ &= [J_{\mathcal{D}\mathbf{f}}(\mathbf{x}) \mid \mathbf{0}] + [J_{\mathcal{D}\mathbf{p}}^{\mathbf{x}}(\mathbf{x}, \varepsilon) \mid J_{\mathcal{D}\mathbf{p}}^{\varepsilon}(\mathbf{x}, \varepsilon)], \end{aligned} \quad (19)$$

where  $J_{\mathcal{D}\mathbf{p}}^{\mathbf{x}}(\mathbf{x}, \varepsilon)$  (resp.  $J_{\mathcal{D}\mathbf{p}}^{\varepsilon}(\mathbf{x}, \varepsilon)$ ) is the Jacobian matrix of  $\mathcal{D}\mathbf{p}$  with respect to  $\mathbf{x}$  (resp.  $\varepsilon$ ). By construction the Jacobian matrix  $J_{\mathcal{D}\mathbf{p}}^{\varepsilon}(\mathbf{x}, \varepsilon)$  of the system  $\mathbf{p} = (\Lambda_i^{\mathbf{x}}(p_j))_{1 \leq i, j \leq \mu}$  is, up to a reordering of the rows and columns, a block diagonal matrix with  $s$  blocks of the form

$$(\Lambda_i^{\mathbf{x}}[\mathbf{b}_j])_{1 \leq i, j \leq \mu}.$$

As  $\mathcal{D}$  is dual to the basis  $\mathbf{b}$ ,  $(A_i^x[\mathbf{b}_j](\zeta, \mathbf{0}))_{1 \leq i, j \leq \mu}$  is the identity matrix, the Jacobian  $J_{\mathcal{D}\mathbf{p}}^\varepsilon(\mathbf{x}, \varepsilon)$  evaluated at  $(\zeta, \mathbf{0})$  is, up to a reordering of the rows and columns, the identity matrix of dimension  $\mu s$ .

Using decomposition (19), we easily deduce the following property:

**Lemma 3.** *The  $\mu s \times \mu s$  Jacobian matrix  $J_{\mathcal{D}\mathbf{g}}^\varepsilon(\mathbf{x}, \varepsilon)$  is of full rank  $\mu s$  at  $(\zeta, \mathbf{0})$ .*

Another interesting property is the following [15]:

**Lemma 4.** *The  $\mu s \times n$  Jacobian matrices  $J_{\mathcal{D}\mathbf{g}}^x(\mathbf{x}, \varepsilon)$  and  $J_{\mathcal{D}\mathbf{f}}^x(\mathbf{x}, \varepsilon)$  are of full rank  $n$  at  $(\zeta, \mathbf{0})$ .*

We are going to use these properties to construct sub-systems of  $\mathcal{D}\mathbf{g}$  with a simple root “above”  $\zeta$ .

The columns of  $J_{\mathcal{D}\mathbf{g}}(\mathbf{x}, \varepsilon)$  are indexed by the variables  $(\mathbf{x}, \varepsilon)$ , while the rows are indexed by the polynomials  $g_{ik}$ . We construct the following systems:

- (a) Let  $\mathcal{D}\mathbf{f}^I$  be a subsystem of  $\mathcal{D}\mathbf{f}$  s.t. the corresponding  $n$  rows of  $J_{\mathcal{D}\mathbf{f}}(\zeta)$  are linearly independent (Lemma 4 implies that such a subset exists). We denote by  $I = \{(i_1, k_1), \dots, (i_n, k_n)\}$  their indices.
- (b) Let  $\mathcal{D}\tilde{\mathbf{g}}(\mathbf{x}, \tilde{\varepsilon})$  be the square system formed by removing the variables  $\varepsilon_{k_1, i_1}, \dots, \varepsilon_{k_n, i_n}$  from  $\mathcal{D}\mathbf{g}(\mathbf{x}, \varepsilon)$ . Therefore the Jacobian  $J_{\mathcal{D}\tilde{\mathbf{g}}}(\mathbf{x}, \tilde{\varepsilon})$  derives from  $J_{\mathcal{D}\mathbf{g}}(\mathbf{x}, \varepsilon)$ , after purging the columns indexed by  $\varepsilon_{k_1, i_1}, \dots, \varepsilon_{k_n, i_n}$  and it's  $(i_j, k_j)$ -th row becomes  $[\nabla(A_{i_j}^x \tilde{g}_{i_j, k_j})^T \mid \mathbf{0}]$ .

A first consequence is the following result, giving a  $n \times n$  system deduce from the initial system  $\mathbf{f}$ , with a simple root at  $\zeta$ :

**Theorem 5 (Deflation Theorem 1 [15]).** *Let  $\mathbf{f}(\mathbf{x})$  be a  $n$ -variate polynomial system with an  $\mu$ -fold isolated zero at  $\mathbf{x} = \zeta$ . Then the  $n \times n$  system  $\mathcal{D}\mathbf{f}^I(\mathbf{x}) = 0$ , defined in (a), has a simple root at  $\mathbf{x} = \zeta$ .*

*Example 9.* In our running example, we expand the rectangular Jacobian matrix of 6 polynomials in  $(x_1, x_2)$ . Choosing the rows corresponding to  $f_1$  and  $(d_1 - d_2^2 - d_1 d_2 - d_1^2)[f_1]$ , we find a non-singular minor, hence the resulting system  $(f_1, 2x_1)$  has a regular root at  $\zeta = (0, 0)$ .

The deflated system  $\mathcal{D}\mathbf{f}^I(\mathbf{x}) = 0$  is a square system in  $n$  variables. Contrarily to the deflation approach in [12, 4], we do not introduce new variables and one step of deflation is provably sufficient. The trade-off is that here we assume that exact dual elements are pointed at by indices  $I$ , so as to be able to compute the original multiple root with high accuracy.

On the other hand, when the coefficients are machine numbers, an exact multiple root is unlikely to exist. In the following theorem, we introduce new variables that will allow us later to derive an approximate deflation method. The need to introduce new variables comes from the fact that in practice the exact root is not available, or even worse, the input coefficients contain small error. Therefore, our method shall seek for a slightly perturbed system with

an exact multiple zero within a controlled neighborhood of the input, that fits as close as possible to the approximate multiplicity structure of the input system and point.

**Theorem 6 (Deflation Theorem 2 [15]).** *Let  $\mathbf{f}(\mathbf{x})$  be a  $n$ -variate polynomial system with a  $\mu$ -fold isolated root at  $\mathbf{x} = \boldsymbol{\zeta}$ . The square system  $\mathcal{D}\tilde{\mathbf{g}}(\mathbf{x}, \tilde{\mathbf{e}}) = 0$ , as defined in (b), has a regular isolated root at  $(\mathbf{x}, \tilde{\mathbf{e}}) = (\boldsymbol{\zeta}, \mathbf{0})$ .*

Nevertheless, this deflation does differ from the deflation strategy in [12, 4]. There, new variables are added that correspond to coefficients of differential elements, thus introducing a perturbation in the dual basis. This is suitable for exact equations, but, in case of perturbed data, the equations do not actually define a true singular point.

*Example 10.* Consider the system [13] of 3 equations in 2 variables  $f_1 = x_1^3 + x_1x_2^2$ ,  $f_2 = x_1x_2^2 + x_2^3$ ,  $f_3 = x_1^2x_2 + x_1x_2^2$ , and the singular point  $(0, 0)$  of multiplicity equal to 7.

Suppose that the point is given. Using 3 and 2 we derive the primal-dual pair

$$\mathcal{D} = (1, d_1, d_2, d_1^2, d_1d_2, d_2^2, \underline{d_2^3} + d_1^3 + d_1^2d_2 - d_1d_2^2),$$

where  $\underline{d_2^3}$  is underlined to show that it corresponds to  $x_2^3$  in the primal monomial basis  $\mathcal{B} = (1, x_1, x_2, x_1^2, x_1x_2, x_2^2, x_2^3)$ . The biggest matrix used, in depth 4, was of size  $9 \times 8$ , while Macaulay's method terminates with a matrix of size  $30 \times 15$ .

To deflate the root, we construct the augmented system  $\mathcal{D}\mathbf{f}$  of 21 equations. The  $21 \times 2$  Jacobian matrix  $J_{\mathcal{D}\mathbf{f}}(\mathbf{x})$  is of rank 2 and a full-rank minor consists of the rows 4 and 5. Therefore, we find the system  $(d_1^2[f_1], d_1d_2[f_1]) = (3x_1, 2x_2)$  which deflates  $(0, 0)$ . Note that even though both equations of the deflated system derive from  $f_1$ , the functionals used on  $f_1$  are computed using all initial equations.

The perturbed equations are then

$$\begin{aligned} g_1 &= f_1 + \varepsilon_{1,1} + \varepsilon_{1,2}x_1 + \varepsilon_{1,3}x_2 + \varepsilon_{1,4}x_2^2 + \varepsilon_{1,5}x_2^3 \\ g_2 &= f_2 + \varepsilon_{2,1} + \varepsilon_{2,2}x_1 + \varepsilon_{2,3}x_2 + \varepsilon_{2,4}x_1^2 + \varepsilon_{2,5}x_1x_2 + \varepsilon_{2,6}x_2^2 + \varepsilon_{2,7}x_2^3 \\ g_3 &= f_3 + \varepsilon_{3,1} + \varepsilon_{3,2}x_1 + \varepsilon_{3,3}x_2 + \varepsilon_{3,4}x_1^2 + \varepsilon_{3,5}x_1x_2 + \varepsilon_{3,6}x_2^2 + \varepsilon_{3,7}x_2^3 \end{aligned}$$

and the resulting system  $\mathcal{D}\mathbf{g}$  has a simple root at  $(\boldsymbol{\zeta}, \mathbf{0})$ .

## 5 Approximate multiple point

In real-life applications it is common to work with approximate inputs. Also, there is the need to (numerically) decide if an (approximate) system possesses a single (real) root in a given domain, notably for use in subdivision-based algorithms, e.g. [19, 16].

In the regular case, Smale's  $\alpha$ -theory, extending Newton's method, can be used to answer this problem, also partially extended to singular cases in [7], using zero clustering. Another option is to use the following certification test, based on the verification method of Rump [22, Th. 2.1]:

**Theorem 7 ([10, 22] Krawczyk-Rump Theorem).** *Let  $\mathbf{f} \in R^n$ ,  $R = \mathbb{K}[\mathbf{x}]$ , be a polynomial system and  $\zeta^* \in \mathbb{R}^n$  a real approximate regular isolated point. Given an interval domain  $Z \in \mathbb{IR}^n$  containing  $\zeta^* \in \mathbb{R}^n$ , and an interval matrix  $M \in \mathbb{IR}^{n \times n}$  whose  $i$ -th column  $M_i$  satisfies*

$$\nabla f_i(Z) \subseteq M_i \quad \text{for } i = 1 \dots, n$$

*then the following holds: If the interval domain*

$$V_{\mathbf{f}}(Z, \zeta^*) := -J_{\mathbf{f}}(\zeta^*)^{-1} \mathbf{f}(\zeta^*) + (I - J_{\mathbf{f}}(\zeta^*)^{-1} M)Z \quad (20)$$

*is contained in  $\overset{\circ}{Z}$ , the interior of  $Z$ , then there is a unique  $\zeta \in Z$  with  $\mathbf{f}(\zeta) = \mathbf{0}$  and the Jacobian matrix  $J_{\mathbf{f}}(\zeta) \in M$  is non-singular.*

In our implementation we use this latter approach, since it is suitable for inexact data and suits best with the perturbation which is applied. In particular, it coincides with the numerical scheme of [22] in the univariate case.

In the case of an isolated multiple point of a polynomial system, we applied a deflation to transform it into a regular root of an extended system. The theorem is applied to the system of 6, using an (approximate) structure  $\mathcal{D}$ . The resulting range of the  $\varepsilon$ -parameters encloses a system that attains a single multiple root of that structure. Hence the domain for  $\varepsilon$ -variables reflects the distance of the input system from a precise system with local structure  $\mathcal{D}$ . Therefore, we obtain a perturbed system in a neighborhood of the input together with a numerically controlled bound on the perturbation coefficients, with a unique multiple root having a prescribed multiplicity.

If the multiple point is known approximately, we use implicitly Taylor's expansion of the polynomials at this approximate point to deduce the dual basis, applying the algorithm of the previous section. The following computation can be applied:

- At each step, the solutions of linear system (10, i-iii) are computed via Singular Value Decomposition. Using a given threshold, we determine the numerical rank and an orthogonal basis of the solutions from the last singular values and the last columns of the right factor of the SVD.
- For the computation of the monomials which define the equations (2, iii) at the next step, we apply QR decomposition on the transpose of the basis to extract a non-zero maximal minor. The monomials indexing this minor are used to determine constraints (10, i-iii). A similar numerical technique is employed in [25], for Macaulay's method.

*Example 11.* Let  $f_1 = x_1^2x_2 - x_1x_2^2$ ,  $f_2 = x_1 - x_2^2$ . The verification method of [22] applies a linear perturbation to this system, but fails to certify the root  $\zeta = (0, 0)$ .

We consider an approximate point  $\zeta^* = (.01, .002)$  and we compute the approximate multiplicity structure

$$\mathcal{D} = (A_1, \dots, A_4) = (1.0, 1.0d_2, \underline{1.0d_1} + 1.0d_2^2, \underline{1.0d_1d_2} + 1.0d_2^3).$$

The augmented system  $\mathbf{g}(\mathbf{x}) = (A_j[f_i]) = (f_1, 2.0x_1x_2 - 1.0x_2^2 - 1.0x_1, 2.0x_1 - 2.0x_2, 1.0x_1 - 1.0x_2^2, f_2, -2.0x_2, 0., 0.)$  has a Jacobian matrix:

$$J_{\mathbf{g}}(\zeta^*)^T = \begin{bmatrix} .00 & .016 & -.99 & 2.0 & 1.0 & 0 & 0 & 0 \\ .00 & -.02 & .016 & -2.0 & -.004 & -2.0 & 0 & 0 \end{bmatrix}$$

with a non-zero minor at the third and forth row. Using this information, we apply the following perturbation to the original system:

$$\begin{aligned} g_1 &= x_1^2x_2 - x_1x_2^2 + \varepsilon_{11} + \varepsilon_{12}x_2 \\ g_5 &= x_1 - x_2^2 + \varepsilon_{21} + \varepsilon_{22}x_2 + \varepsilon_{23}x_1 + \varepsilon_{24}x_1x_2 \end{aligned}$$

Thus  $\mathbf{g}(x_1, x_2, \varepsilon_{11}, \varepsilon_{12}, \varepsilon_{21}, \varepsilon_{22}, \varepsilon_{23}, \varepsilon_{24})$ , computed as before, is a square system with additional equations:

$$\begin{aligned} g_2 &= 1.0x_1^2 - 2.0x_1x_2 + 1.0\varepsilon_{12} \\ g_3 &= 2.0x_1x_2 - 1.0x_2^2 - 1.0x_1 \\ g_4 &= 2.0x_1 - 2.0x_2 \\ g_6 &= -2.0x_2 + 1.0\varepsilon_{22} + 1.0x_1\varepsilon_{24} \\ g_7 &= 1.0\varepsilon_{23} + 1.0x_2\varepsilon_{24} \\ g_8 &= 1.0\varepsilon_{24} \end{aligned}$$

Now take the box  $Z_1 = [-.03, .05] \times [-.04, .04] \times [-.01, .01]^6$ . We apply Th. 7 on  $\mathbf{g}$ , i.e. we compute  $V_{\mathbf{g}}(Z_1, \zeta^*)$ . For the variable  $\varepsilon_{21}$  the interval is  $[-.015, .15] \not\subseteq (-.01, .01)$ , therefore we don't get an answer.

We shrink a little  $Z_1$  down to  $Z_2 = [-.03, .05] \times [-.02, .02] \times [-.01, .01]^6$  and we apply again Th. 7, which results in

$$V_{\mathbf{g}}(Z_2, (\zeta^*, \mathbf{0})) = \begin{bmatrix} [-.004, .004] \\ [-.004, .004] \\ [-.001, .001] \\ [-.007, .007] \\ [-.006, .006] \\ [-.009, .009] \\ [-.00045, .00035] \\ [.0, .0] \end{bmatrix} \subseteq \overset{\circ}{Z}_2,$$

thus we certify that the input equations admit a perturbation of magnitude of .01, so that the perturbed system has a unique exact root within the interval  $[-.03, .05] \times [-.02, .02]$ .



## 6 Experimentation

We have implemented the presented algorithms in MAPLE. It can compute (approximate) dual bases by means of Macaulay's method as well as the integration method, and it can derive the augmented system defined in Th. 6. Then Krawczyk-Rump's interval method is used to verify the root.

*Example 12.* Let, as in [11, 13],

$$\begin{aligned} f_1 &= 2x_1 + 2x_1^2 + 2x_2 + 2x_2^2 + x_3^2 - 1, \\ f_2 &= (x_1 + x_2 - x_3 - 1)^3 - x_1^3, \\ f_3 &= (2x_1^3 + 2x_2^2 + 10x_3 + 5x_3^2 + 5)^3 - 1000x_1^5. \end{aligned}$$

The point  $(0, 0, -1)$  occurs with multiplicity equal to 18, in depth 7. The final matrix size with our method is  $54 \times 37$ , while Macaulay's method ends with a  $360 \times 165$  matrix.

If the objective is to deflate as efficiently as possible, then one can go step by step: First compute a basis of  $\mathcal{D}_1$  and stop the process. We get the evaluation 1 and 2 first order functionals, which we apply to  $f_1$ . We arrive at

$$(\mathbf{1}[f_1], (d_2 - d_1)[f_1], (d_1 + d_3)[f_1]) = (f_1, -4x_1 + 4x_2, 2 + 4x_1 + 2x_3)$$

and we check that the Jacobian determinant is 64, thus we have a deflated system only with a partial local structure. The condition number of the Jacobian matrix is also very satisfactory, with a value of around 5.55.

The recent paper [8], implementing the dialytic deflation method produces a deflated system of size  $75 \times 48$  for this instance, with a condition number of order  $10^6$ .

*Example 13.* Consider the equations (taken from [4, DZ3]):

$$\begin{aligned} f_1 &= 14x_1 + 33x_2 - 3\sqrt{5}(x_1^2 + 4x_1x_2 + 4x_2^2 + 2) + \sqrt{7} + x_1^3 + 6x_1^2x_2 + 12x_1x_2^2 + 8x_2^3, \\ f_2 &= 41x_1 - 18x_2 - \sqrt{5} + 8x_1^3 - 12x_1^2x_2 + 6x_1x_2^2 - x_2^3 + 3\sqrt{7}(4x_1x_2 - 4x_1^2 - x_2^2 - 2) \end{aligned}$$

and take an approximate system  $\tilde{\mathbf{f}}$  with those coefficients rounded to 6 digits. A 5-fold zero of  $\tilde{\mathbf{f}}$  rounded to 6 digits is  $\zeta^* = (1.50551, .365278)$ .

Starting with the approximate system and with a tolerance of .001, we compute the basis

$$\begin{aligned} \mathcal{D} = & (1, d_1 + .33d_2, d_1^2 + .33d_1d_2 + .11d_2^2, d_1^3 + .33d_1^2d_2 + .11d_1d_2^2 + .03d_2^3 - \\ & - 1.54d_2, d_1^4 + .33d_1^3d_2 + .11d_1^2d_2^2 + .03d_1d_2^3 + .01d_2^4 - 1.54d_1d_2 - 1.03d_2^2) \end{aligned}$$

having 4 correct digits, with respect to the initial exact system, and the primal counterpart  $\mathcal{B} = (1, x_1, x_1^2, x_1^3, x_1^4)$ .

We form the deflated system (b), with  $I = \{(3, 1), (5, 1)\}$ , i.e. the 3rd and 5th dual element on  $f_1$  have non-null Jacobian. By adding 8 new variables, the system is perturbed as:

$$g_{1,1} = \tilde{f}_1 + \varepsilon_{1,1} + \varepsilon_{1,2}(x_1 - \zeta_1^*) + \varepsilon_{1,4}(x_1 - \zeta_1^*)^3,$$

$$g_{2,1} = \tilde{f}_2 + \sum_{i=1}^5 \varepsilon_{2,i}(x_1 - \zeta_1^*)^{i+1}$$

and their derivation with respect to  $\mathcal{D}$ .

We consider a box  $Z$  with center  $= \zeta^*$  and length  $= .004$  at each side. Also, we allow a range  $E = [-.004, .004]^8$  for the variables  $\tilde{\varepsilon}$ . Applying 7 we get a verified inclusion  $V_{\mathbf{g}}(Z \times E, (\zeta^*, \mathbf{0}))$  inside  $Z \times E$  and we deduce that a unique specialization  $\tilde{\varepsilon} \in E$  “fits” the approximate system  $\tilde{\mathbf{f}}$  to the multiplicity structure  $\mathcal{D}$ .

Indeed, one iteration of Newton’s method on  $\mathbf{g}(\mathbf{x}, \varepsilon)$  gives the approximate point  $\zeta = (1.505535473, .365266196)$  and corresponding values for  $\varepsilon_0 \in E$ , such that  $\zeta$  is a 9–digit approximation of the multiple root of the perturbed system  $\mathbf{g}(\mathbf{x}, \varepsilon_0)$ .

In Table 1 we run dual basis computation on the benchmark set of [4]. Multiplicity, matrix sizes at termination step and computation time is reported. One sees that there is at least an order of gain in the running time using the primal-dual approach.

System	$\mu/n$	MM’11		Mourrain’97		Macaulay	
cmbs1	11/3	27 × 23	.18s	27 × 33	.95s	105 × 56	1.55s
cmbs2	8/3	21 × 17	.08s	21 × 24	.39s	60 × 35	.48s
mth191	4/3	10 × 9	.03s	10 × 12	.07s	30 × 20	.14s
decker2	4/2	5 × 5	.02s	5 × 8	.05s	20 × 15	.10s
Ojika2	2/3	6 × 5	.02s	6 × 6	.03s	12 × 10	.04s
Ojika3	4/3	12 × 9	.07s	12 × 12	.27s	60 × 35	.59s
KSS	16/5	155 × 65	8.59s	155 × 80	40.41s	630 × 252	70.03s
Capr.	4/4	22 × 13	.28s	22 × 16	.47s	60 × 35	2.34s
Cyclic-9	4/9	104 × 33	1.04s	104 × 36	5.47s	495 × 220	31.40s
DZ1	131/4	700 × 394	14m	700 × 524	26m	4004 × 1365	220m
DZ2	16/3	43 × 33	.68s	43 × 48	4.38s	360 × 165	25.72s
DZ3	5/2	6 × 6	.04s	6 × 10	.23s	30 × 21	.79s

**Table 1.** Benchmark systems from [3], reporting matrix size at the last step of computing a dual local basis, and overall time for primal-dual computation. The computations are done using Maple. Observe that Macaulay’s method results in a matrix of size  $n \binom{p-1+n}{p-1} \times \binom{p+n}{p}$ , in contrast to a matrix of size  $(\frac{n(n-1)}{2} \mu + n) \times \mu(n-1) + 1$  for the primal-dual approach.

**Acknowledgements.** This research has received funding from the EU’s 7<sup>th</sup> Framework Programme [FP7/2007-2013], Marie Curie Initial Training Network SAGA, grant n° [PITN-GA-2008-214584].

## References

1. L. Alberti, B. Mourrain, and J. Wintz. Topology and arrangement computation of semi-algebraic planar curves. *Comput. Aided Geom. Des.*, 25:631–651, November 2008.
2. M.F. Atiyah and I.G. MacDonald. *Introduction to Commutative Algebra*. Addison-Wesley, 1969.
3. B. H. Dayton, T.-Y. Li, and Z. Zeng. Multiple zeros of nonlinear systems. *Math. Comput.*, 80(276):2143–2168, 2011.
4. B. H. Dayton and Z. Zeng. Computing the multiplicity structure in solving polynomial systems. In *ISSAC '05: Proceedings of the 2005 International Symposium on Symbolic and Algebraic Computation*, pages 116–123, New York, NY, USA, 2005. ACM.
5. M. Elkadi and B. Mourrain. *Introduction à la résolution des systèmes d'équations algébriques*, volume 59 of *Mathématiques et Applications*. Springer-Verlag, 2007.
6. W. J. Gilbert. Newton's method for multiple roots. *Computers & Graphics*, 18(2):227–229, 1994.
7. M. Giusti, G. Lecerf, B. Salvy, and J.-C. Yakoubsohn. On location and approximation of clusters of zeros: Case of embedding dimension one. *Foundations of Computational Mathematics*, 7:1–58, 2007. 10.1007/s10208-004-0159-5.
8. W. Hao, A. J. Sommese, and Z. Zeng. Algorithm 931: An algorithm and software for computing multiplicity structures at zeros of nonlinear systems. *ACM Trans. Math. Softw.*, 40(1):5:1–5:16, October 2013.
9. L. V. Kantorovich. Functional analysis and applied mathematics. *Uspekhi Matematicheskikh Nauk*, 3(6):89–185, 1948.
10. R. Krawczyk. Newton-algorithmen zur bestimmung von nullstellen mit fehlerschranken. *Computing*, 4(3):187–201, 1969.
11. G. Lecerf. Quadratic newton iteration for systems with multiplicity. *Foundations of Computational Mathematics*, 2:247–293, 2002.
12. A. Leykin, J. Verschelde, and Zhao A. Newton's method with deflation for isolated singularities of polynomial systems. *Theoretical Computer Science*, 359(1-3):111 – 122, 2006.
13. A. Leykin, J. Verschelde, and A. Zhao. Higher-order deflation for polynomial systems with isolated singular solutions. In A. Dickenstein, F.-O. Schreyer, and A.J. Sommese, editors, *Algorithms in Algebraic Geometry*, volume 146 of *The IMA Volumes in Mathematics and its Applications*, pages 79–97. Springer New York, 2008.
14. F.S. Macaulay. *The algebraic theory of modular systems*. Cambridge Univ. Press, 1916.
15. A. Mantzaflaris and B. Mourrain. Deflation and Certified Isolation of Singular Zeros of Polynomial Systems. In A. Leykin, editor, *International Symposium on Symbolic and Algebraic Computation (ISSAC)*, pages 249–256, San Jose, CA, United States, June 2011. ACM New York.
16. A. Mantzaflaris, B. Mourrain, and E. Tsigaridas. Continued fraction expansion of real roots of polynomial systems. In *Proceedings of the 2009 Conference on Symbolic-Numeric Computation, SNC '09*, pages 85–94, New York, NY, USA, 2009. ACM.

17. M. G. Marinari, T. Mora, and H.M. Möller. Gröbner duality and multiplicities in polynomial system solving. In *Proceedings of the 1995 International Symposium on Symbolic and Algebraic Computation*, ISSAC '95, pages 167–179, New York, NY, USA, 1995. ACM.
18. B. Mourrain. Isolated points, duality and residues. *Journal of Pure and Applied Algebra*, 117-118:469 – 493, 1997.
19. B. Mourrain and J. P. Pavone. Subdivision methods for solving polynomial equations. *J. Symb. Comput.*, 44:292–306, March 2009.
20. T. Ojika, S. Watanabe, and T. Mitsui. Deflation algorithm for the multiple roots of a system of nonlinear equations. *Journal of Mathematical Analysis and Applications*, 96(2):463 – 479, 1983.
21. S.R. Pope and A. Szanto. Nearest multivariate system with given root multiplicities. *Journal of Symbolic Computation*, 44(6):606 – 625, 2009.
22. S. Rump and S. Graillat. Verified error bounds for multiple roots of systems of nonlinear equations. *Numerical Algorithms*, 54:359–377, 2010. 10.1007/s11075-009-9339-3.
23. H. J. Stetter. Analysis of zero clusters in multivariate polynomial systems. In *Proceedings of the 1996 International Symposium on Symbolic and Algebraic Computation*, ISSAC '96, pages 127–136, New York, NY, USA, 1996. ACM.
24. X. Wu and L. Zhi. Computing the multiplicity structure from geometric involutive form. In *Proceedings of the twenty-first International Symposium on Symbolic and Algebraic Computation*, ISSAC '08, pages 325–332, New York, NY, USA, 2008. ACM.
25. Z. Zeng. The closedness subspace method for computing the multiplicity structure of a polynomial system. In D. Bates, G. Besana, S. Di Rocco, and C. Wampler, editors, *Interactions of Classical and Numerical Algebraic Geometry*, volume 496 of *Contemporary Mathematics*, pages 347–362. Am. Math. Society, Providence, RI, 2009.